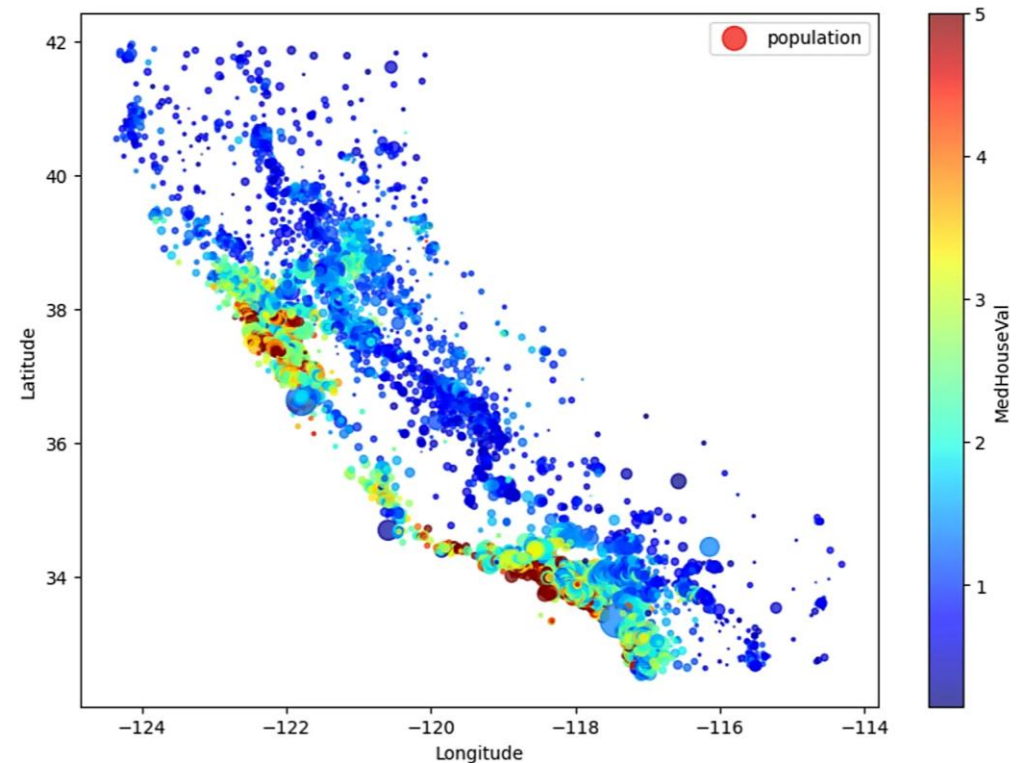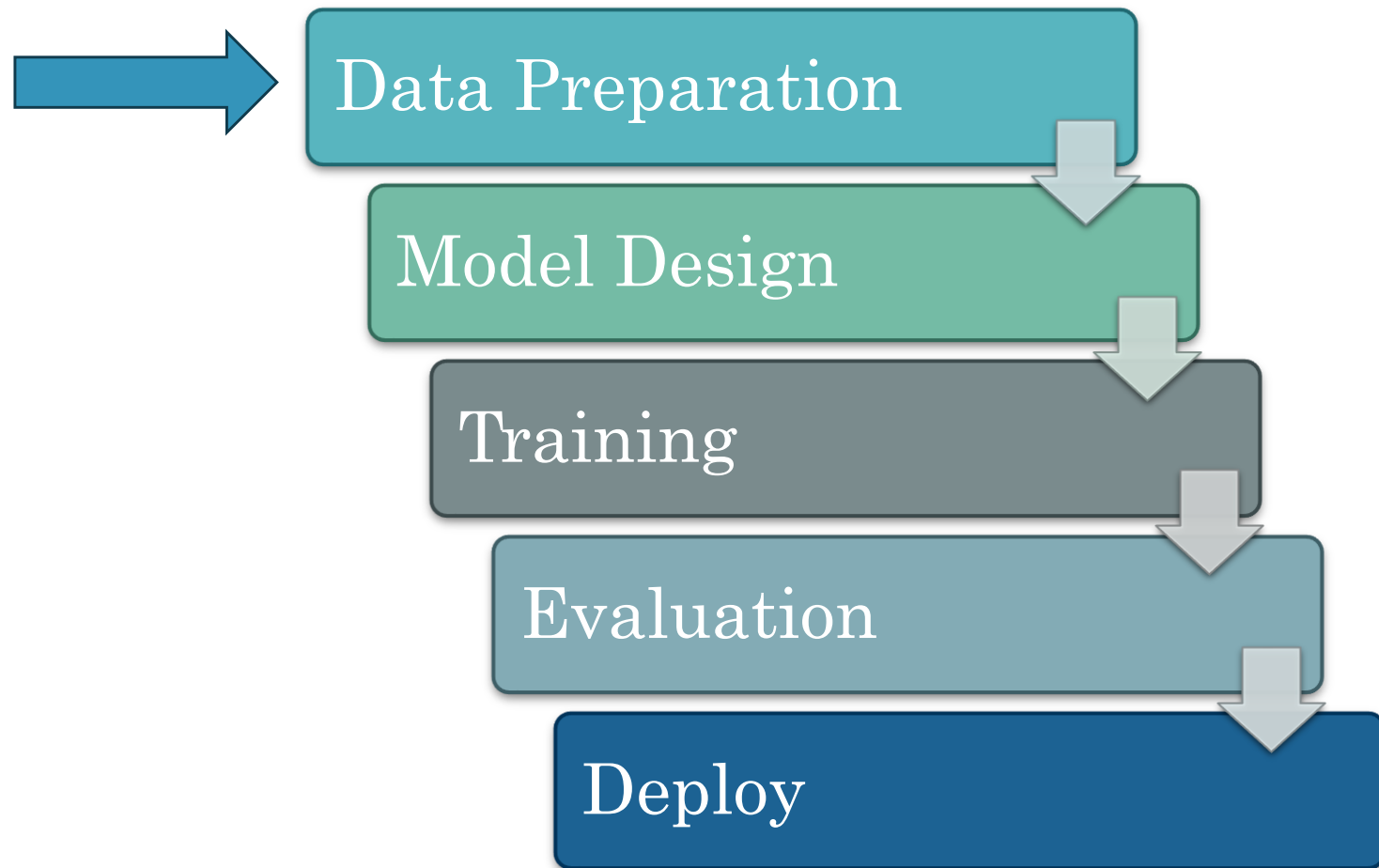# End-to-end Project

## Machine Learning Course

Saeed Mohagheghi

# Project

- Predict **House Price** in California (in 1990!)

- Dataset: **California Housing** [download link]
  - Source:         1990 U.S. Census data
  - Instances:    20,640 districts
  - **Features**:    8 numeric / 1 categorical:
    - Latitude, Longitude
    - Median housing age
    - Total rooms, Total bedrooms
    - Population
    - Households
    - Median income
    - Ocean proximity (categorical)
  - **Target**:         Median house value



→ Book's Chapter 2

# ML Pipeline

# Data Preparation

**Data Analysis**

- Structure and Attributes
- Visualization (Geographical)
- Look for Correlations
- Attribute Combinations
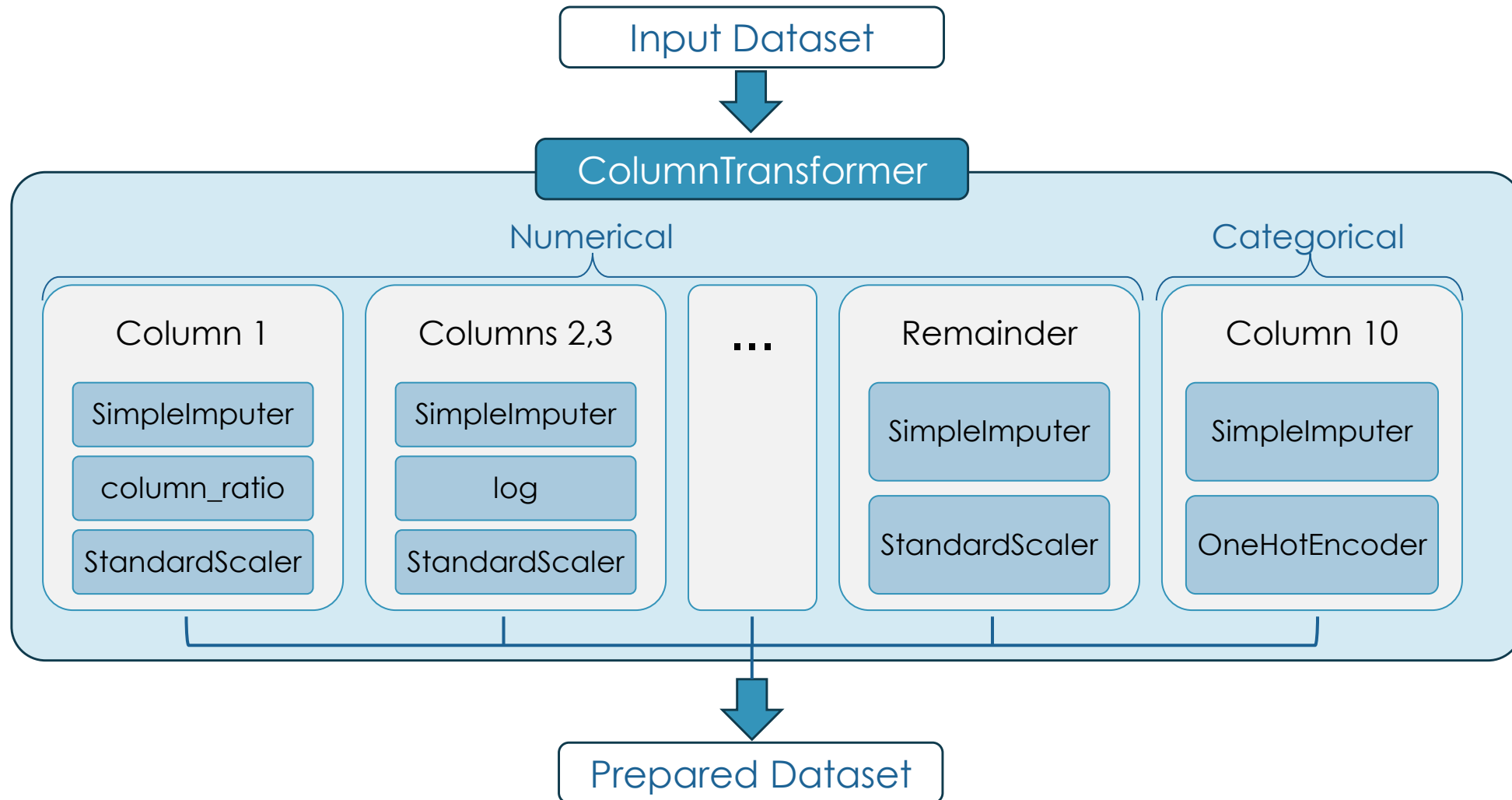
**Prepare for ML**

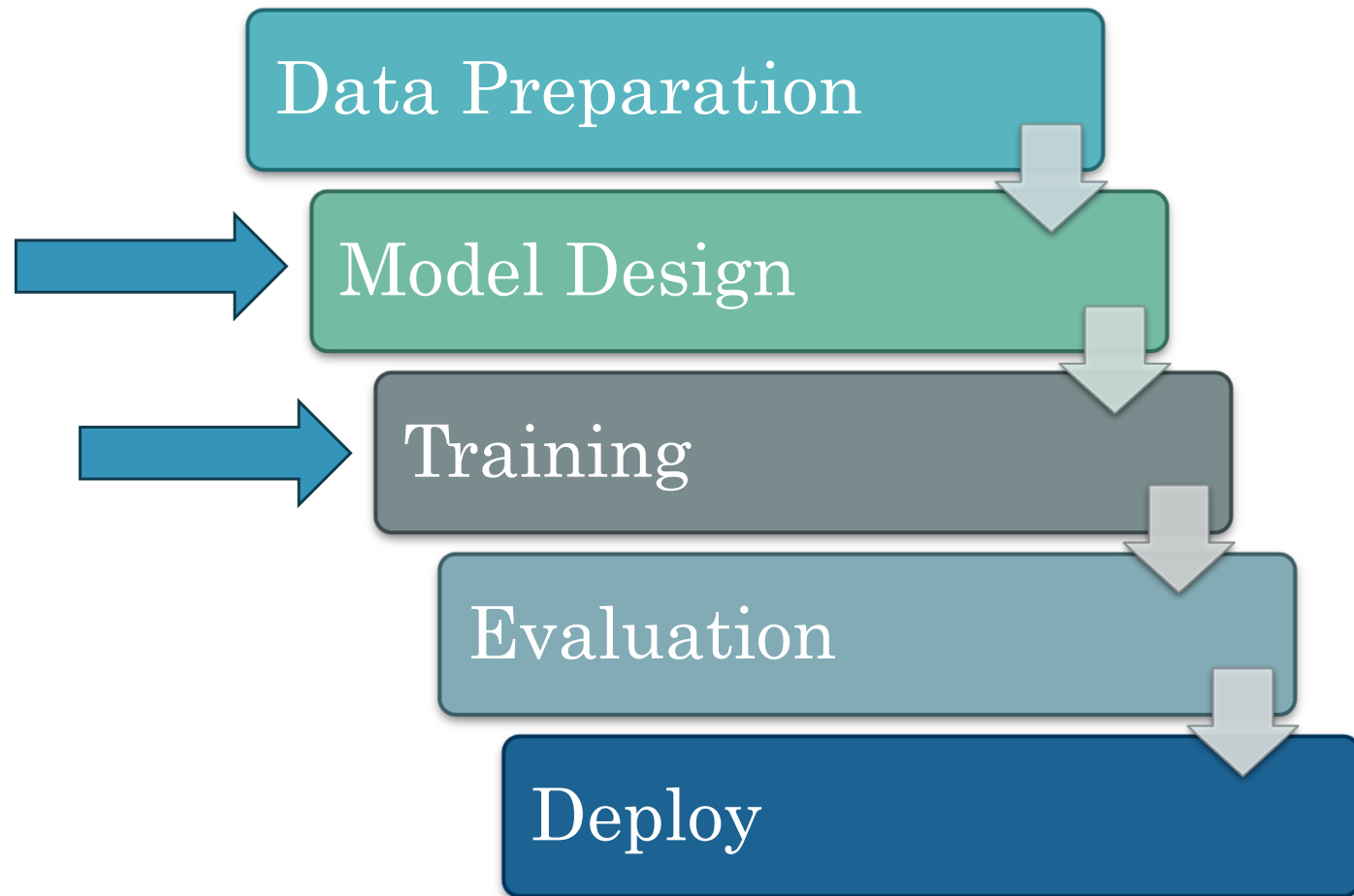- Cleaning
- Train / Test Splits
- Handle Categorical Data
- Feature Scaling and Transformation

4

# Data Preparation Pipeline

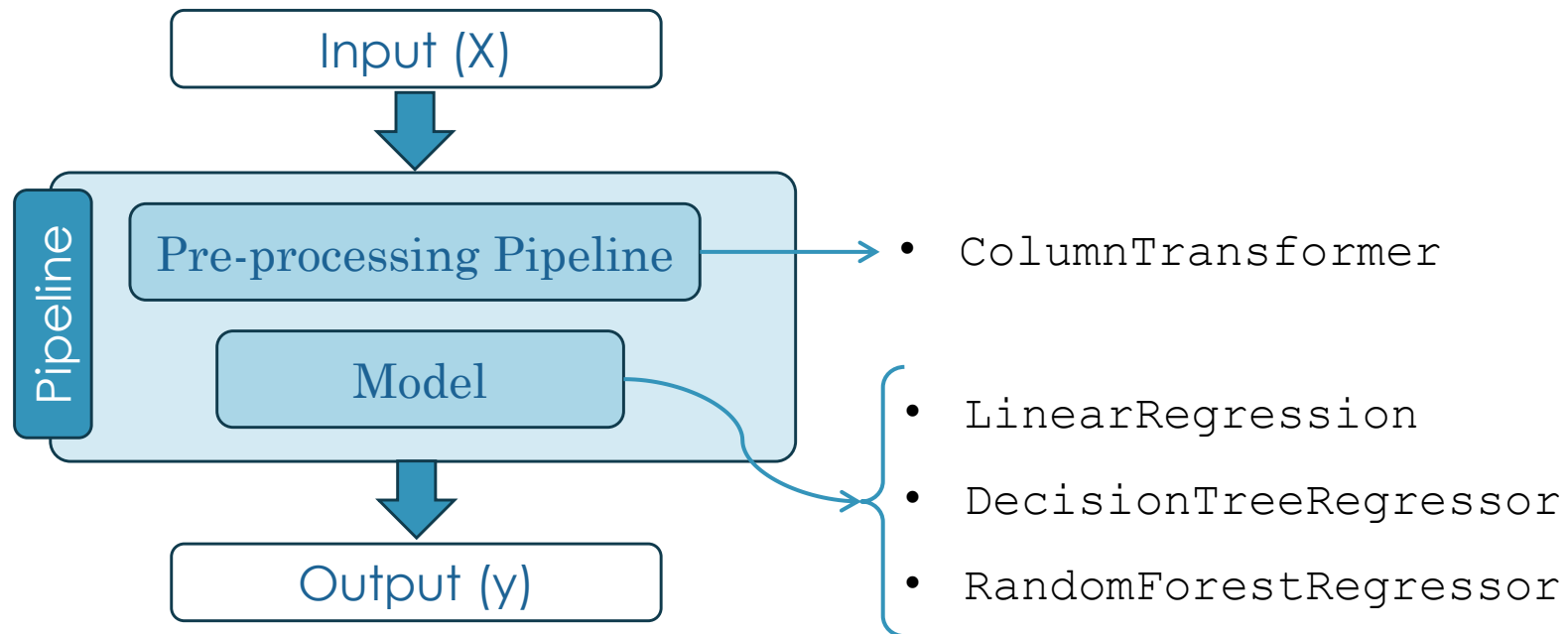# ML Pipeline

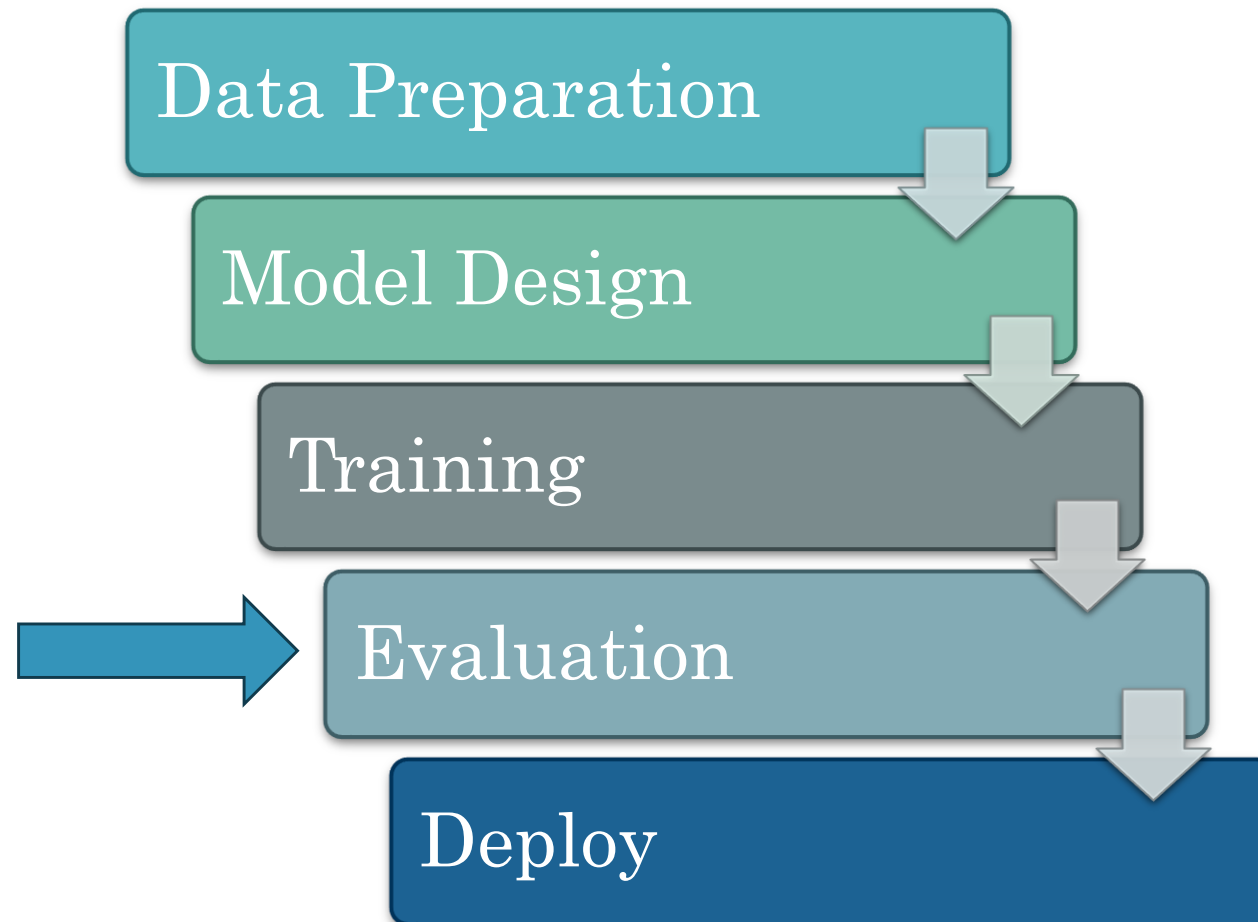# Full Pipeline

- Full Pipeline
  - to be used in Training / Evaluation / Deployment

# ML Pipeline

# Evaluation

- Using <u>Metrics</u>                               (Evaluate on Train/Val/Test Subset Separately)

  - 
    ```python
    from sklearn.metrics import root_mean_squared_error

    predicted_labels = pipeline.predict(housing)

    rmse = root_mean_squared_error( true_labels, predicted_labels )
    ```

- Using <u>K-fold Cross-validation</u>     (Multiple Evaluations on Nonoverlapping Subsets)

  - 
    ```python
    from sklearn.model_selection import cross_val_score

    rmses = -cross_val_score( pipeline, input_data, true_labels,

                             scoring="neg_root_mean_squared_error",

                             cv=k )
    ```